



**University of  
Zurich<sup>UZH</sup>**

**Zurich Open Repository and  
Archive**

University of Zurich  
University Library  
Strickhofstrasse 39  
CH-8057 Zurich  
[www.zora.uzh.ch](http://www.zora.uzh.ch)

---

Year: 2018

---

## **Improving OCR quality of Historical Newspapers with Handwritten Text Recognition Models**

Clematide, Simon ; Ströbel, Phillip

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-167932>

Conference or Workshop Item

Accepted Version



The following work is licensed under a Creative Commons: Attribution 1.0 Generic (CC BY 1.0) License.

Originally published at:

Clematide, Simon; Ströbel, Phillip (2018). Improving OCR quality of Historical Newspapers with Handwritten Text Recognition Models. In: DARIAH-CH Workshop, Neuchâtel, 29 November 2018 - 30 November 2018, University of Neuchâtel.



# Overall impresso pipeline

Main objective: enabling critical text mining to search, extract, process, link, and explore data from print media archives via a unified web interface

## INVOLVES

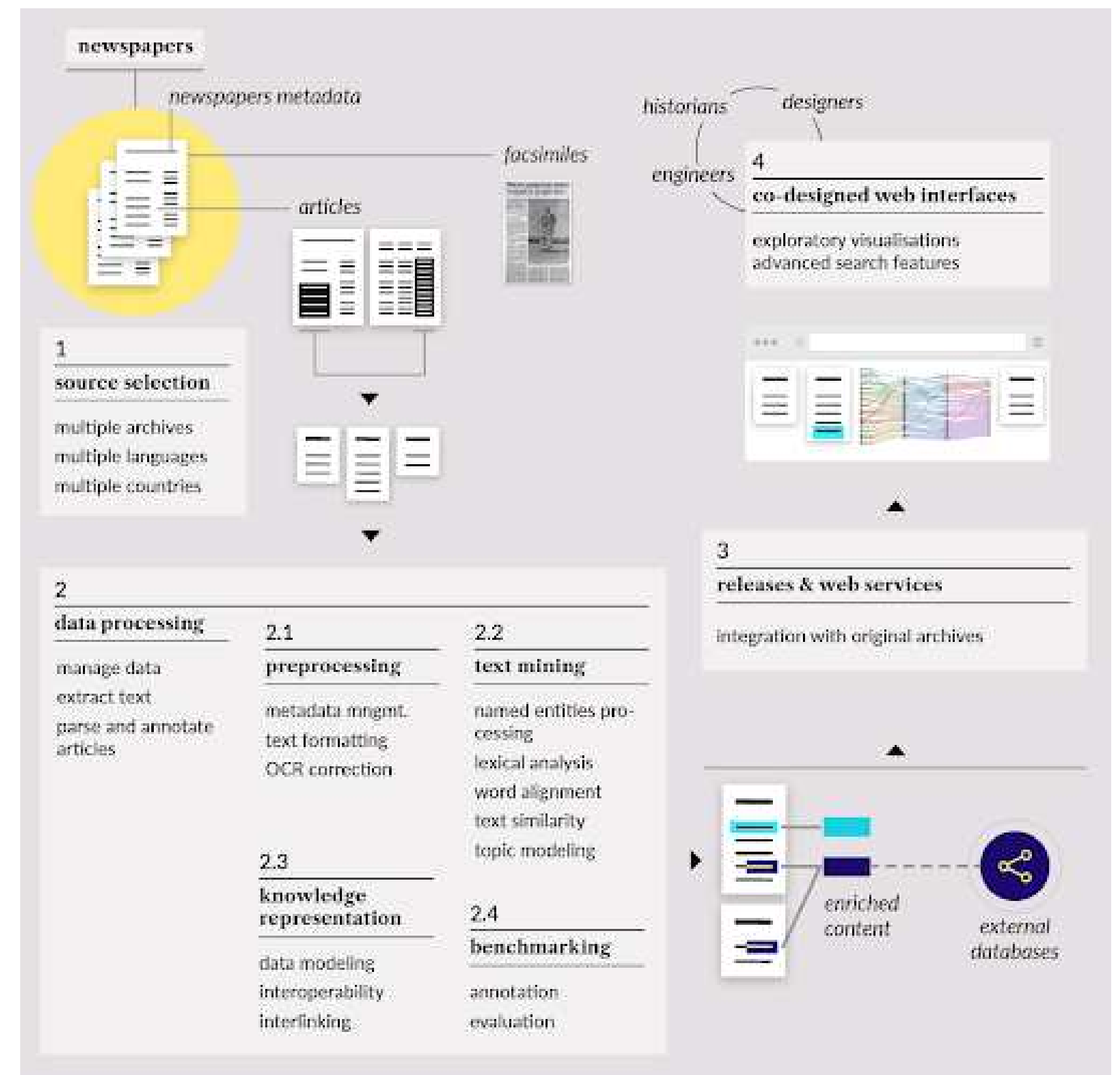
- » computational linguists
- » digital humanists
- » historians
- » web designers

## EXPECTED OUTCOMES

- » natural language processing (NLP) tools dedicated to historical print media written in French and German
- » visualization interfaces for active and goal-oriented exploration and critical analysis of newspaper corpora
- » an application of digital history research on resistance to European integration.

## THEIR COMMON GOAL

- » tackling the challenges of content enrichment and data representation, visualization and analysis, completed by methodological and epistemological reflections



# Surprising success: Researchers report astonishing OCR results on historical newspapers!

## OCR QUALITY OF HISTORICAL NEWSPAPERS



»! , ' ' ' , , ! ^ WWW ' ütchM A ! ' ' i : . . Mex  
 in lAt MroandIM «edrck»»» ; ' ! « ' , ! H ,  
 ! , ' , GroßBrit« . « ! . . chOWtdF . Wm«M« .  
 Hz ; ntshcl« ! , lIn ! i i u u » Nim Norden .  
 Königes pou Isq . M « . » , « ! A A , « . » « « A » ,  
 Prand » u . , Debr« in . . Fl« nkl ; lch (Dü-  
 niour« letzter Uuf« ! l . Edle Vnft« ! l  
 lu« « « « ch . . Gchwell WMatssM , ' G r Z  
 lIdZödn « « « t« vtb« . . « iB , . » w« ch «  
 HZ l , « « ! « i l ; « n . ' 5 weIgs / ich EN .  
 » sellen« ! u Mtlie . habt / Iso« Hob Inre«  
 « i« ! » . . » , , ' die Kapirin bewiest«  
 » apfet . eit ihm unsre Newunderun« « wer-  
 ben h , « l , » und wege« des , bittet , stu-  
 traMjen Maetben« die « i m » . Vestztü .  
 mid bezeugen . - es ; « 3 t « llt . A1M1the-  
 nahck« in Zw . Äbl« gefiel prMM il H .  
 « r« l » ar l mIöM . M . » . « A . Absteßed . « Mri«  
 « Mai« ? b7« A - » ; « i« ge« ? d« ; fi« u«  
 . z . MM« , « « O« d« n « iNluft , erb . demÄl-  
 malWgtt . ihn « u fuch / yu neumen« » u«  
 « i lockmt i« « stiltt iu , wenn Auch  
 « i«cht bh« AWhwint / boHM / SAA . stin t« d« P-  
 tirtel ; S« h« , - Willia« . H« ; ! A ; » ;

## TACKLING OCR WITH HTR TOOLS

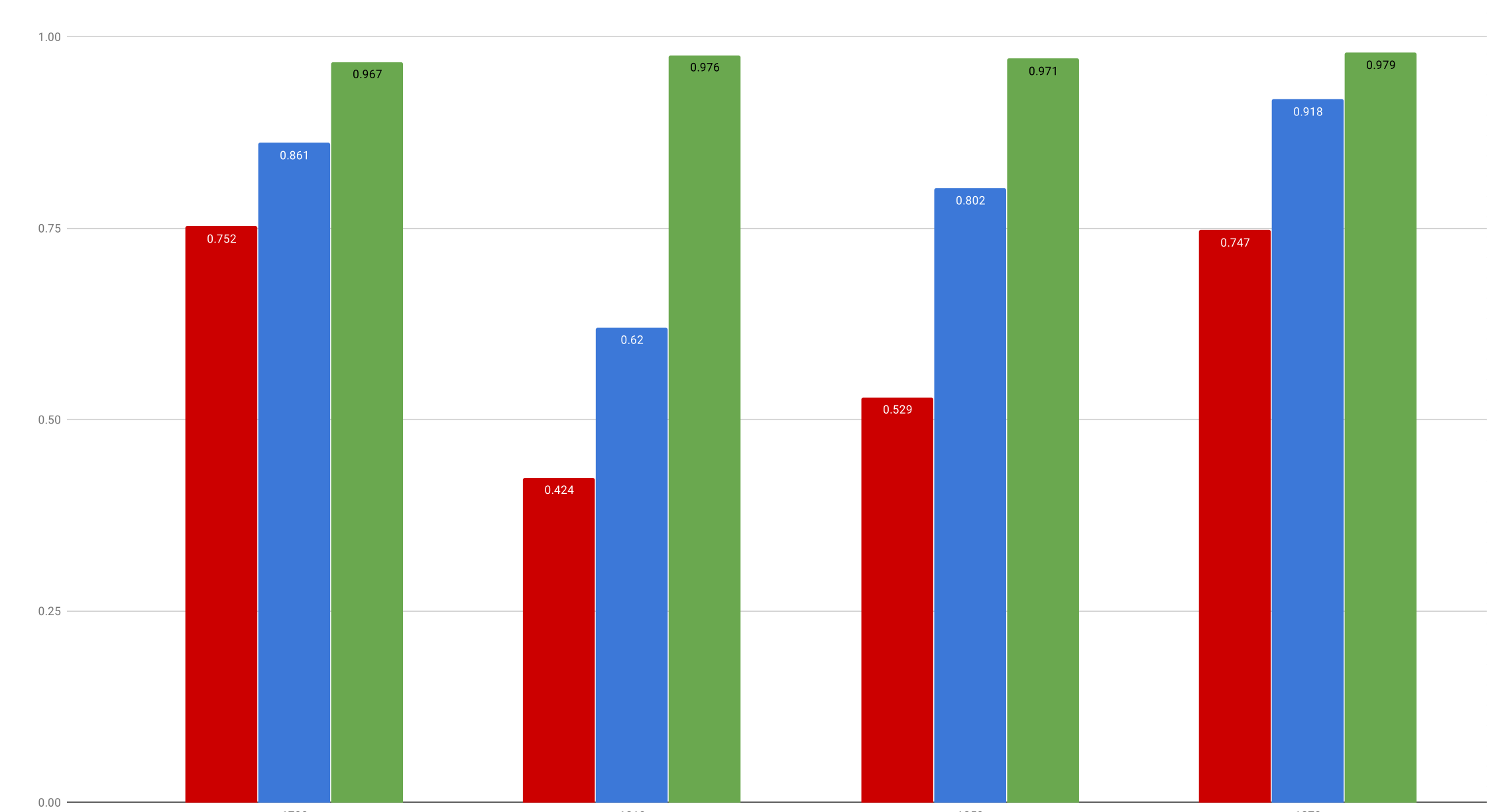


- » 167 front pages from the Neue Zürcher Zeitung (NZZ)
- » we used Transkribus<sup>1</sup> to create a gold standard
- » manual correction of words and baselines
- » training of Handwritten Text
- » Recognition (HTR) model within
- » Transkribus with 158 pages

- 1 <https://www.transkribus.eu>  
2 <https://www.abbyy.com>

## EVALUATION AND COMPARISON

- » bag-of-words f-measure evaluation with TextEval 1.4
  - » [www.primaresearch.org/tools/PerformanceEvaluation](http://www.primaresearch.org/tools/PerformanceEvaluation)
- » compare three different outputs
  - » original OCR by NZZ (**ocr-2005**)
  - » re-OCRised material using ABBYY Finereader<sup>2</sup> (**ocr-2017**)
  - » Transkribus' HTR model (**htr-2018**)



## DISCUSSION & OUTLOOK

- » HTR models significantly increase OCR quality
  - » relatively small gold standards for training purposes suffice for decent OCR
  - » better OCR is beneficial for text mining techniques
- Open questions:
- » Do the HTR models trained on the NZZ perform equally well on other newspapers?
  - » How does occasionally occurring text in antiqua affect OCR quality?

**SUPERVISORS** Frédéric Kaplan, EPFL - Andreas Fickers, C2DH - Martin Volk, UZH

COLLABORATORS Thijs van Beek - Estelle Bunout, PhD - Daniele Guido, M.Sc. - Matteo Romanello, PhD - Paul Schroeder - Phillip Ströbel, M.A.